# ROOT I/O
# Performance and Parallelism

G. Amadio, P. Canal, D. Piparo
for the ROOT Team
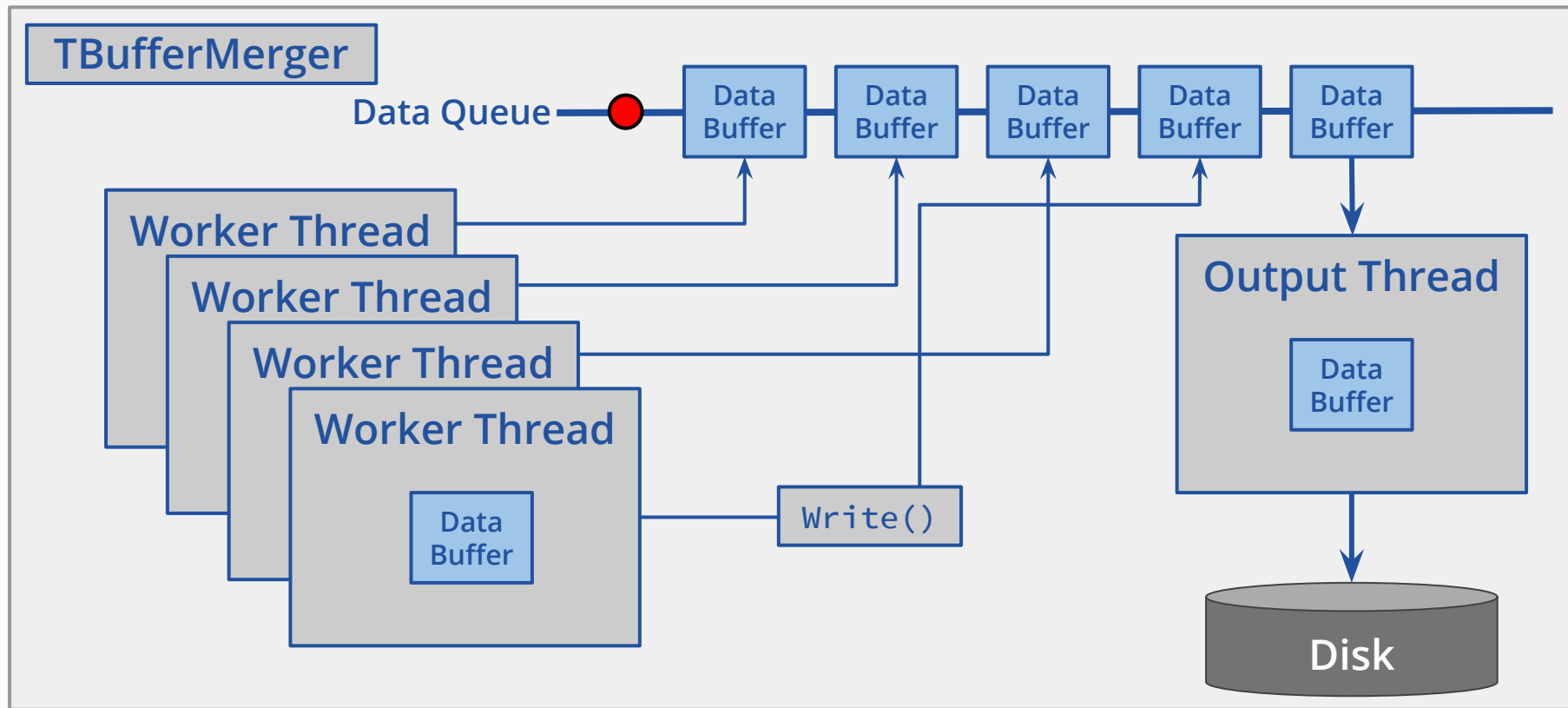
# ROOT
Data Analysis Framework

https://root.cern

Updates coming with ROOT 6.12:

▶ `TBufferMerger` optimizations

▶ ROOT I/O performance improvements

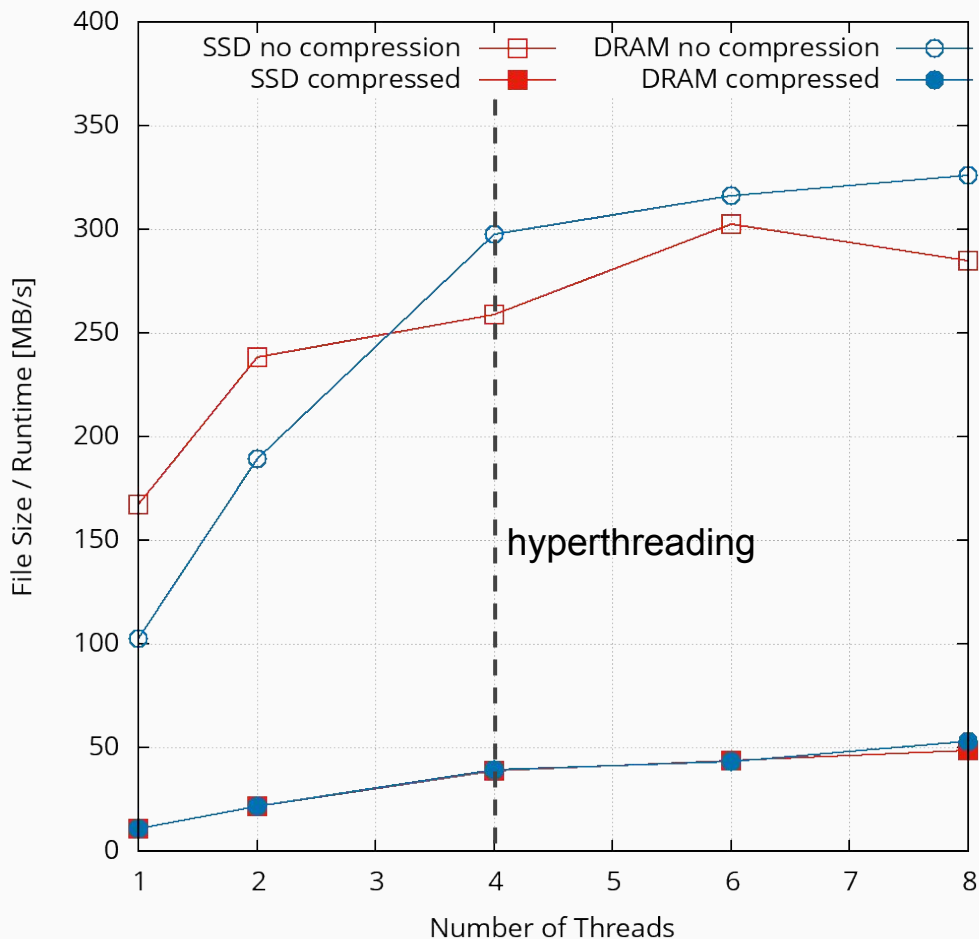- Concurrency of ROOT I/O

- Optimization of `TTree::Fill()`

# Benchmark: TBufferMerger with Random Data

- ▶ Fill a tree with one branch with random numbers

- ▶ Synthetic benchmark that exacerbates the role of I/O by doing only lighweight computations

- ▶ Create ~1GB of data and write out to different media (SSD and DRAM)

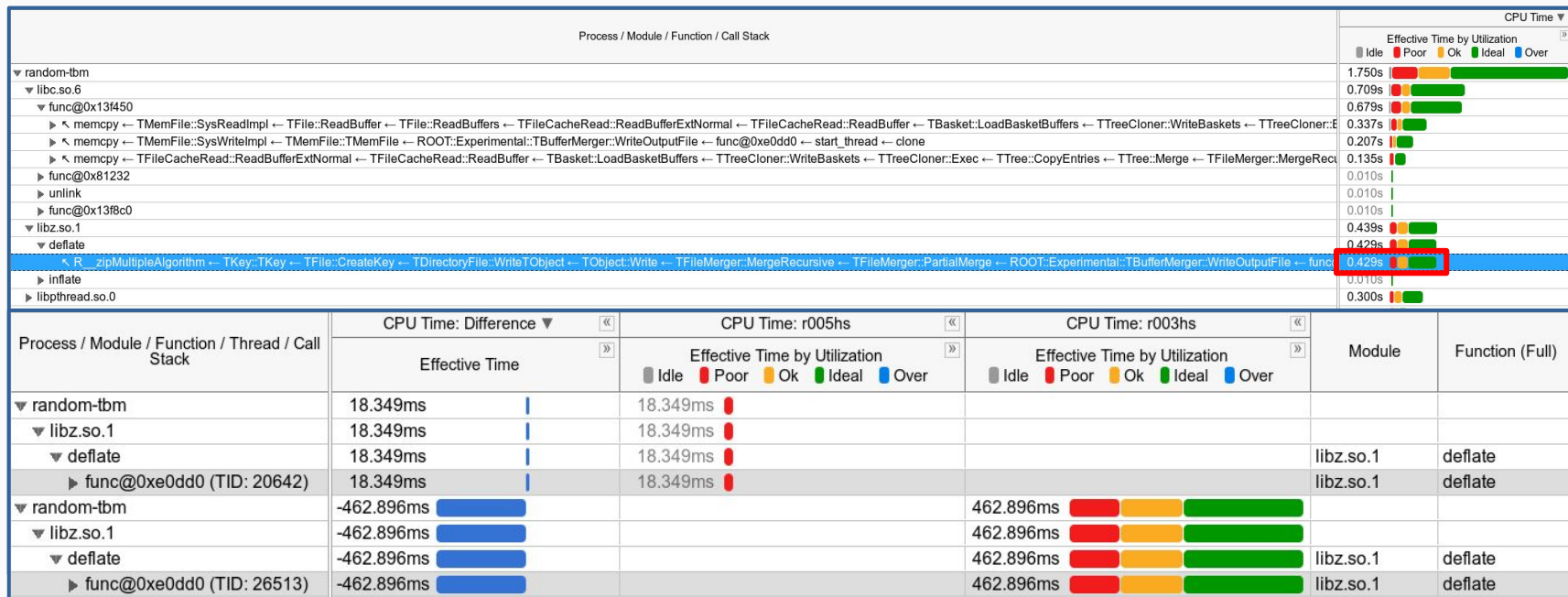- ▶ Quad core laptop Intel® Core i7 4710HQ (2.5GHz, 6M cache)

# TBufferMerger Optimizations

▶ Output thread doing lots of work (compression)
- Solution: add a setting for controlling auto-save point
- Avoids creating too many `TTree` headers, which require compression

▶ Need a way to control queue size or rate-limit
- Solution: add a non-blocking callback mechanism to `TBufferMerger`
- Lets user decide when to create more data-producing tasks by registering a function that gets called everytime a buffer is removed from the merging queue
- Add functions that lets user query the size of the queue
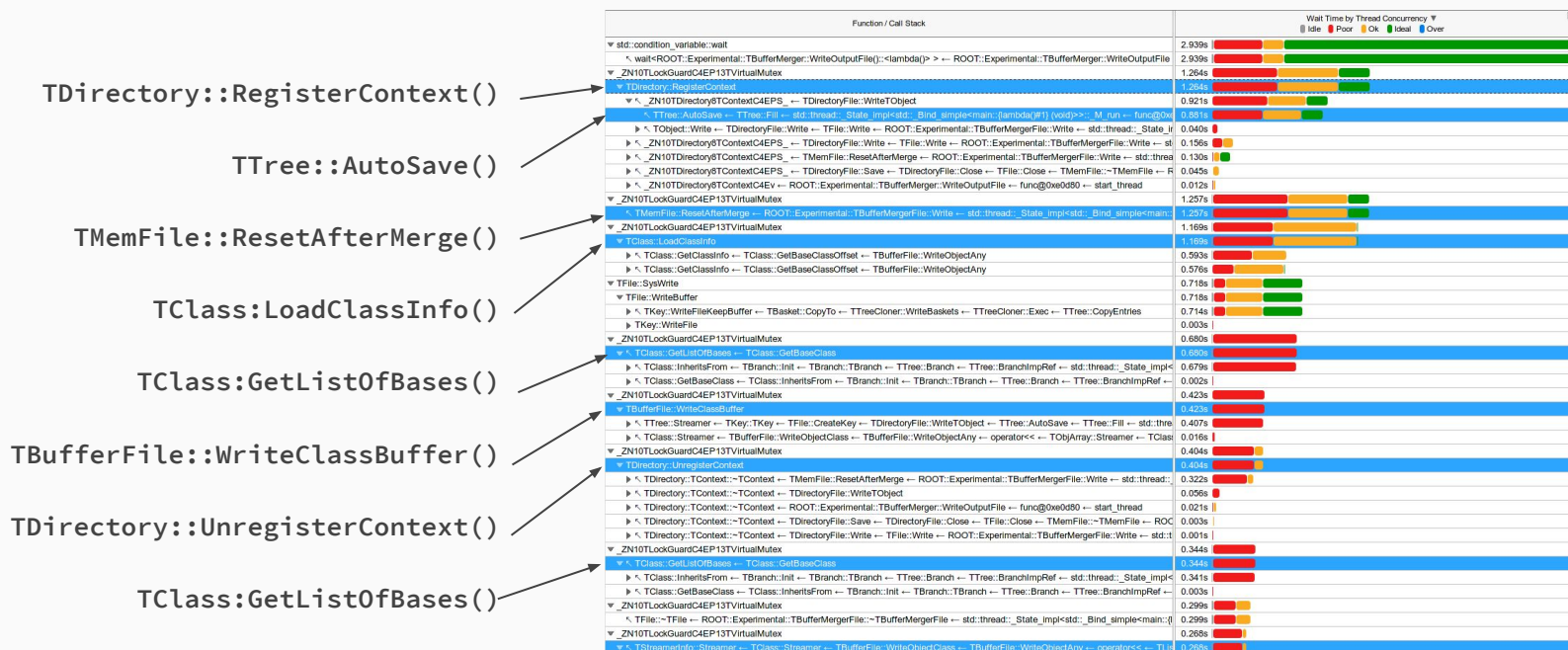- Auto-save forces flush to disk, which also avoids increases in memory

# Output thread compression with small and large auto-save

## Before: many queries to type system → many useless waits



TDirectory::RegisterContext()

TTree::AutoSave()

TMemFile::ResetAfterMerge()

TClass:LoadClassInfo()

TClass:GetListOfBases()

TBufferFile::WriteClassBuffer()

TDirectory::UnregisterContext()

TClass:GetListOfBases()

ROOT-8871

ROOT-9002

7

# Improving the Performance of ROOT I/O

▸ Use simple case with TBufferMerger to optimize ROOT I/O

▸ Same random number generation from before

▸ Reduce number of mutex locks acquired when checking the type system

▸ Reduced from a few hundred locks to a single lock per thread



I/O to disk

before optimization

parallel data generation
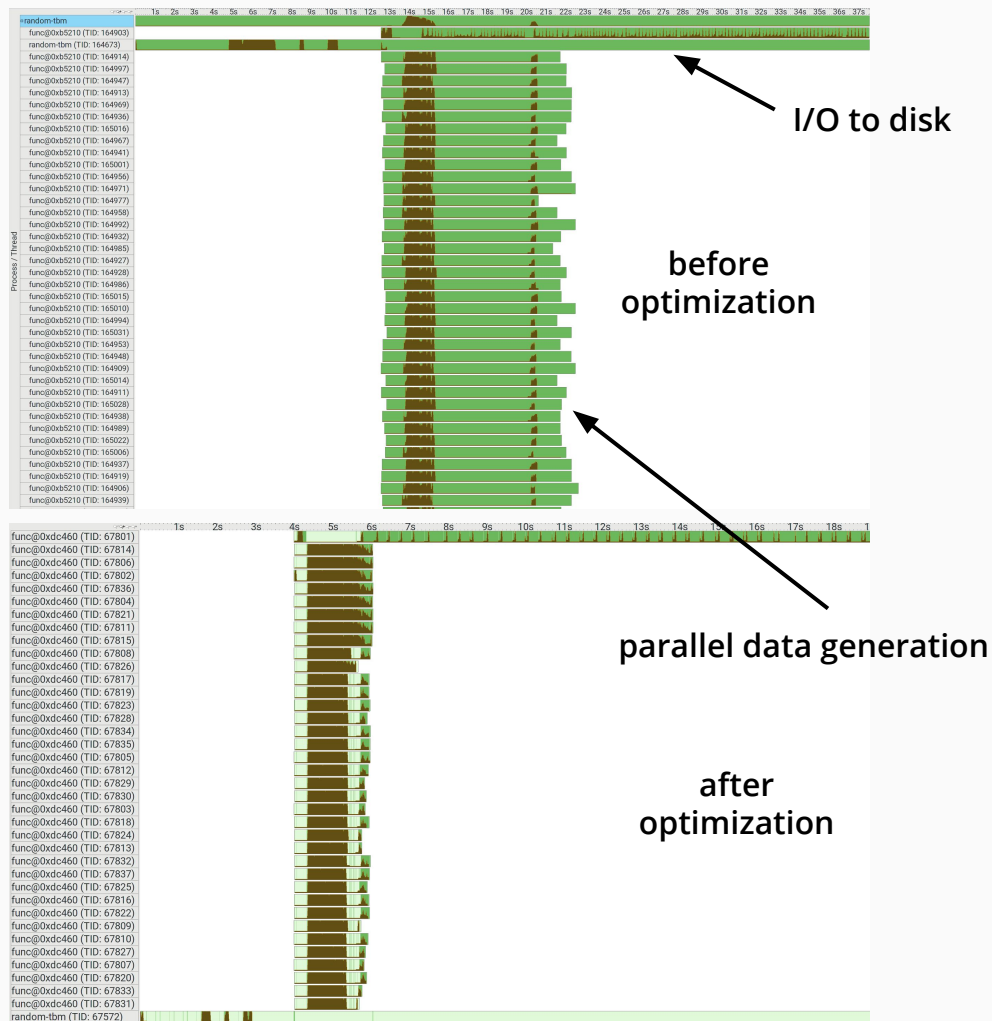
after optimization

# Improving the Performance of ROOT I/O

▸ Use simple case with TBufferMerger to optimize ROOT I/O

▸ Same random number generation from before

▸ Reduce number of mutex locks acquired when checking the type system

▸ Reduced from a few hundred locks to a single lock per thread

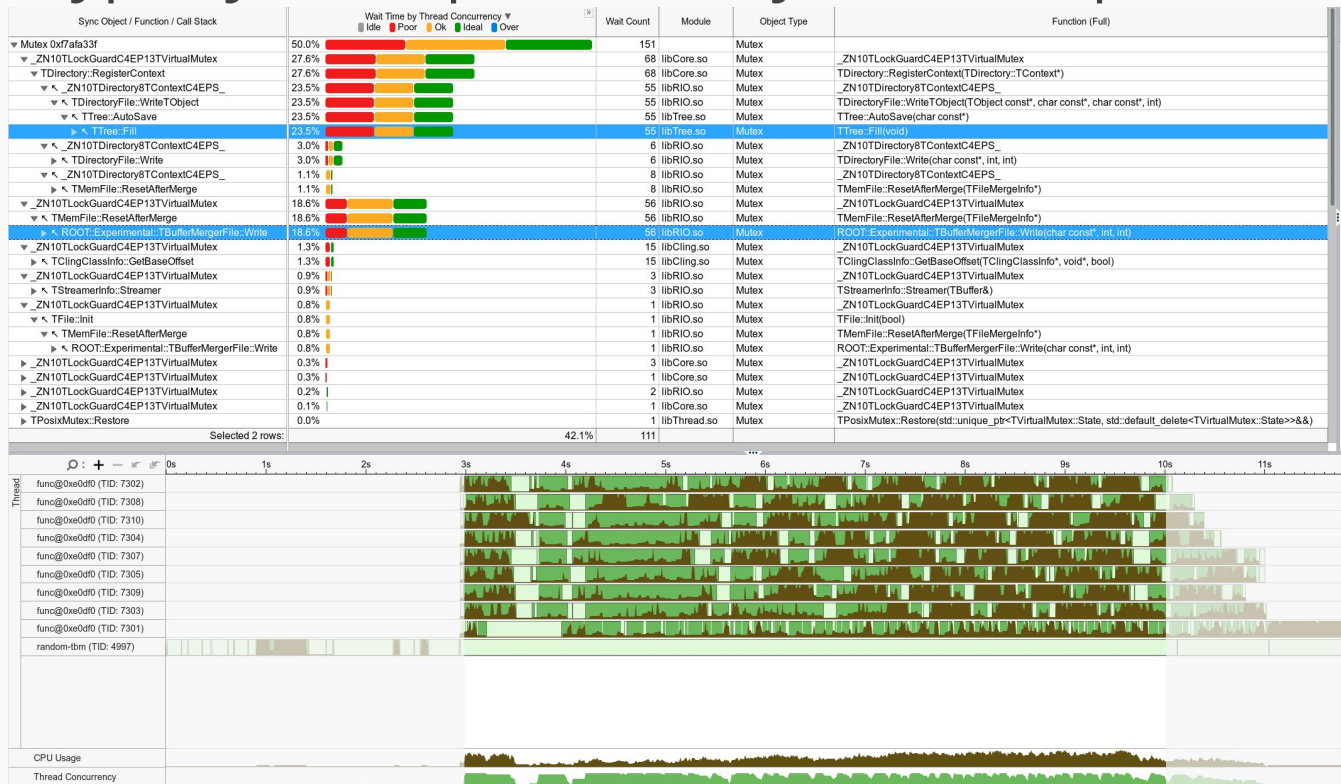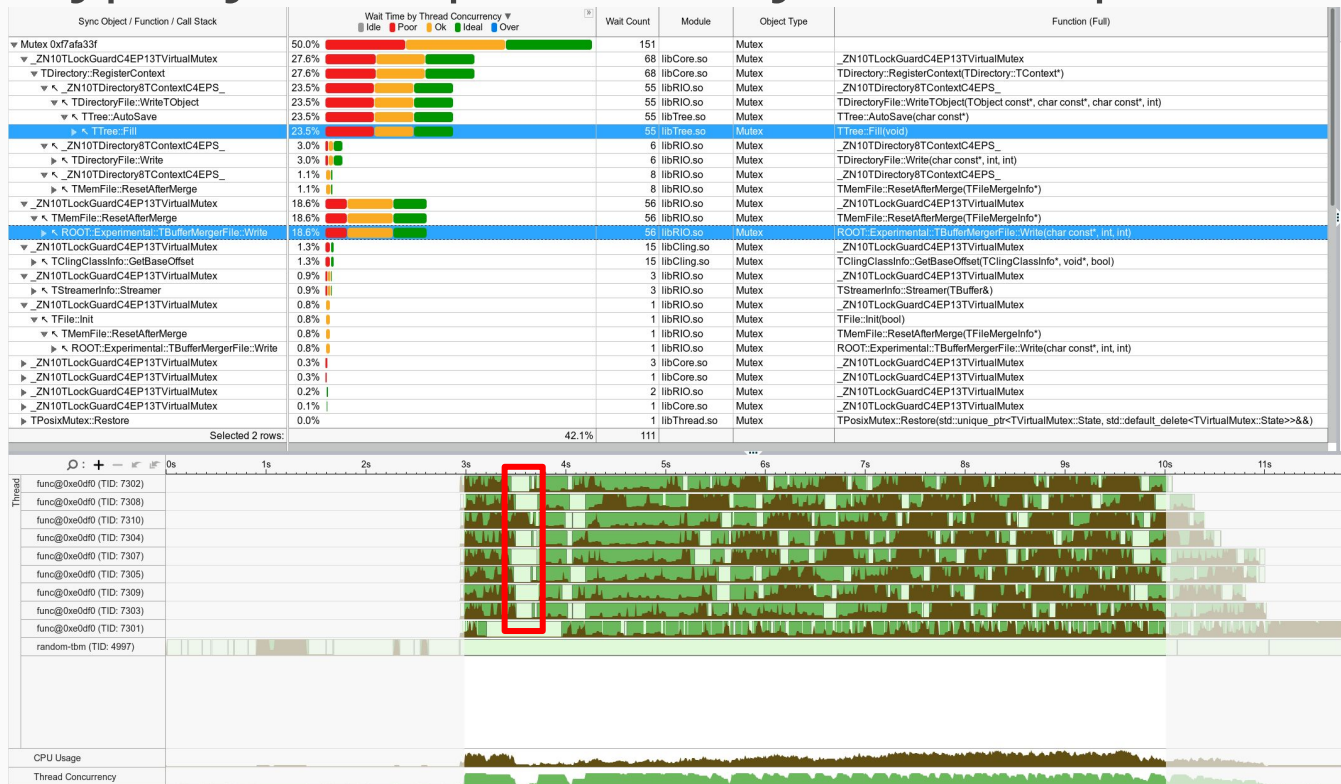**Targeting ROOT 6.12**



I/O to disk

before optimization

parallel data generation

after optimization

No more type system queries, only one wait per thread

No more type system queries, only one wait per thread



11

Divisions take up many CPU cycles for useless work

| Source Line | Source | CPU Time: Total |
| --- | --- | --- |
| | | Effective Time by Utilization |
| | | Idle · Poor · Ok · Ideal · Over |
| 4,446 | ++fEntries; | 1.1% |
| 4,447 | if (fEntries > fMaxEntries) { | 0.1% |
| 4,448 | KeepCircular(); | |
| 4,449 | } | |
| 4,450 | if (gDebug > 0) Info("TTree::Fill", " - A:  %d %lld %lld %lld %lld %lld %lld \n", | 0.3% |
| 4,451 | nbytes, fEntries, fAutoFlush,fAutoSave,GetZipBytes(),fFlushedBytes,fSavedBytes); | |
| 4,452 | | |
| 4,453 | if (fAutoFlush != 0 || fAutoSave != 0) { | 0.2% |
| 4,454 | // Is it time to flush or autosave baskets? | |
| 4,455 | if (fFlushedBytes == 0) { | 0.4% |
| 4,456 | // Decision can be based initially either on the number of bytes | |
| 4,457 | // or the number of entries written. | |
| 4,458 | Long64_t zipBytes = GetZipBytes(); | 0.7% |
| 4,459 | if ((fAutoFlush<0 && zipBytes > -fAutoFlush)  || | 0.2% |
| 4,460 | (fAutoSave <0 && zipBytes > -fAutoSave )  || | 0.1% |
| 4,461 | (fAutoFlush>0 && fEntries%TMath::Max((Long64_t)1,fAutoFlush) == 0) || | 8.9% |
| 4,462 | (fAutoSave >0 && fEntries%TMath::Max((Long64_t)1,fAutoSave)  == 0) ) { | 6.8% |
| 4,463 | | |
| 4,464 | //First call FlushBasket to make sure that fTotBytes is up to date. | |
| 4,465 | FlushBaskets(); | |
| 4,466 | OptimizeBaskets(GetTotBytes(),1,""); | |
| 4,467 | if (gDebug > 0) Info("TTree::Fill","OptimizeBaskets called at entry %lld, fZipBytes=%lld, fFlushedBytes=%lld\n",fEntries,GetZipBytes(),fFlushedBytes); | |
| 4,468 | fFlushedBytes = GetZipBytes(); | |
| 4,469 | fAutoFlush   = fEntries;  // Use test on entries rather than bytes | |
| 4,470 | | |
| 4,471 | // subsequently in run | |
| 4,472 | if (fAutoSave < 0) { | |
| 4,473 | // Set fAutoSave to the largest integer multiple of | |
| 4,474 | // fAutoFlush events such that fAutoSave*fFlushedBytes | |
| 4,475 | // < (minus the input value of fAutoSave) | |
| 4,476 | Long64_t totBytes = GetTotBytes(); | |
| 4,477 | if (zipBytes != 0) { | |
| 4,478 | fAutoSave =  TMath::Max( fAutoFlush, fEntries*((-fAutoSave/zipBytes)/fEntries)); | |
| 4,479 | } else if (totBytes != 0) { | |
| 4,480 | fAutoSave =  TMath::Max( fAutoFlush, fEntries*((-fAutoSave/totBytes)/fEntries)); | |
| 4,481 | } else { | |
| 4,482 | TBufferFile b(TBuffer::kWrite, 10000); | |
| 4,483 | TTree::Class()->WriteBuffer(b, (TTree*) this); | |
| 4,484 | Long64_t total = b.Length(); | |
| 4,485 | fAutoSave =  TMath::Max( fAutoFlush, fEntries*((-fAutoSave/total)/fEntries)); | |
| 4,486 | } | |
| 4,487 | } else if(fAutoSave > 0) { | |
| 4,488 | fAutoSave = fAutoFlush*(fAutoSave/fAutoFlush); | |

12

- ROOT continues to parallelise its I/O subsystem
  - Focus not only on experiments' data processing, but also on analysis
- Parallel writing to single output file via **TBufferMerger**
  - Leveraged by **TDataFrame** already with snapshot action
  - Good performance, can saturate an SSD

Changes already in master for ROOT 6.12 release:

- Optimised **TTree::Fill()** function avoids divisions
- Optimised parallel merging with **TBufferMerger**
- Output thread no longer does excessive compression work
- Callback function allows seamless framework integration

Questions?